

b-jet Identification in the D0 Experiment

Sébastien Greder

Institut Pluridisciplinaire Hubert Curien, Strasbourg

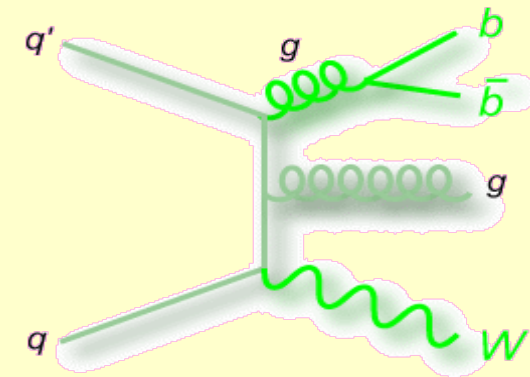
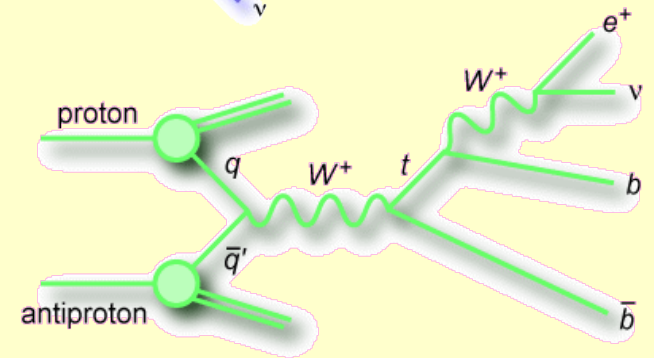
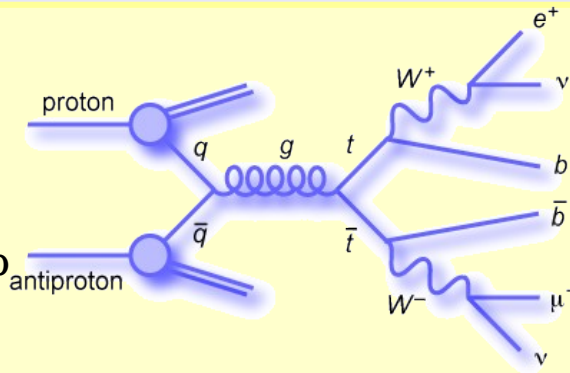
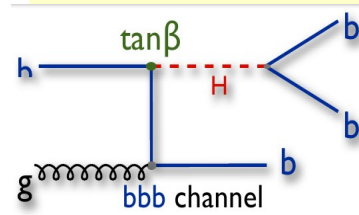
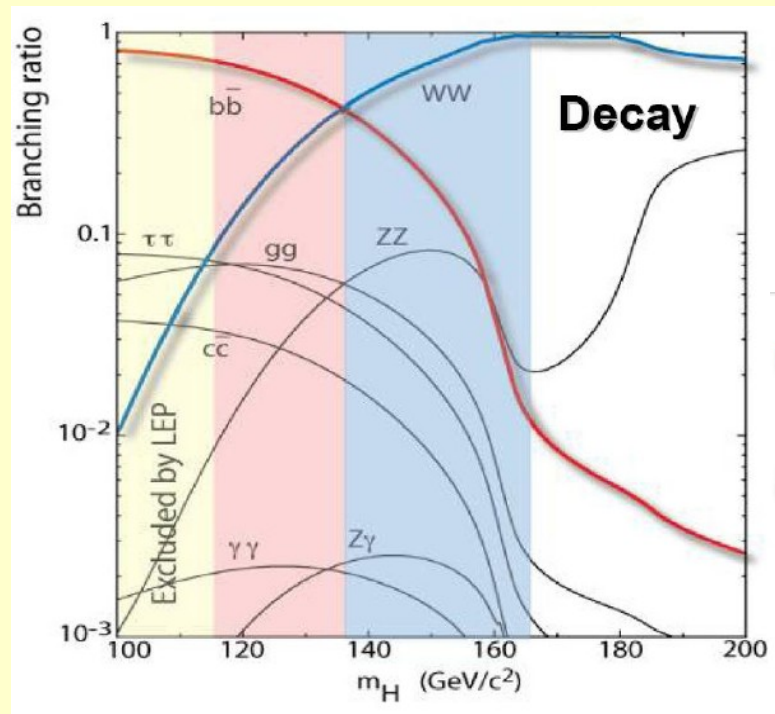
Outline

- Introduction
- The DØ detector
- Algorithms
- Performance measurement
 - B-jet efficiency
 - Fake tag rate
- Conclusion

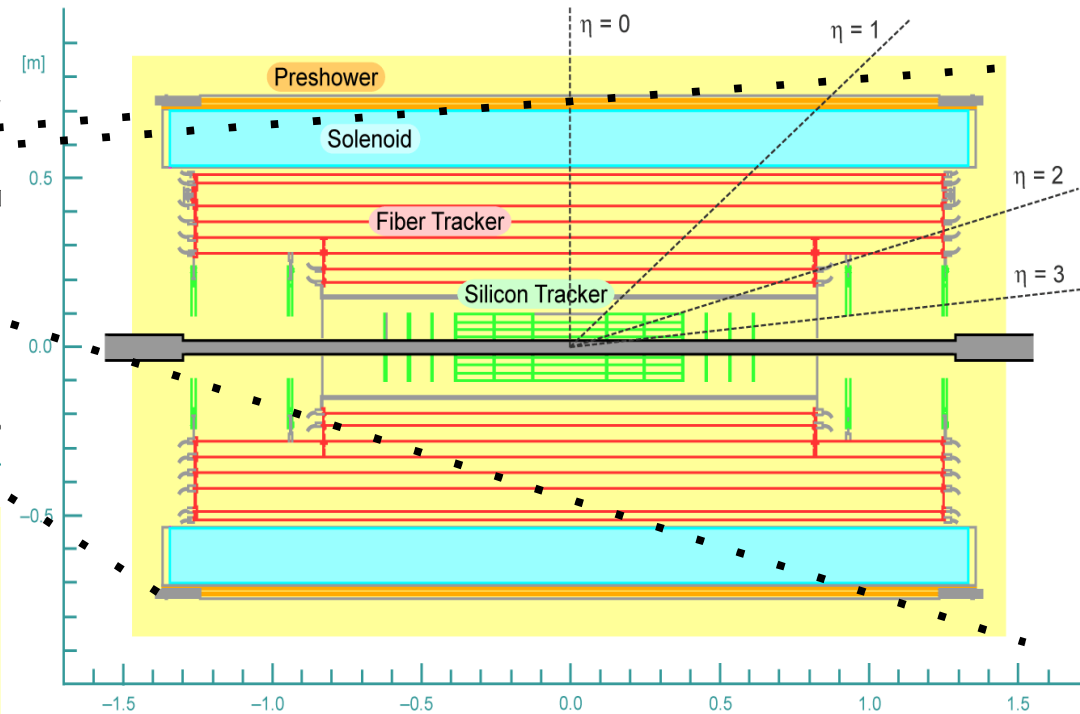
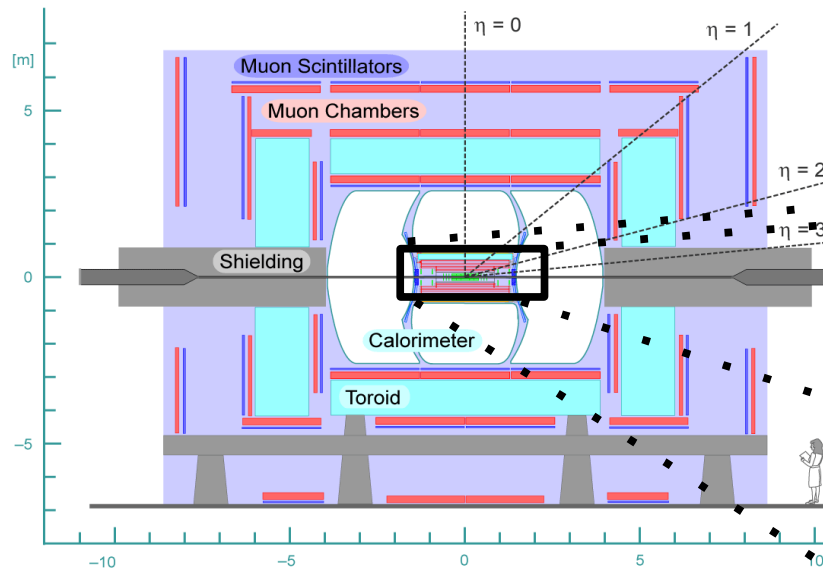
Introduction

Physics

- Top physics: x-section, mass, single-top
- “Backgrounds”: W/Z+heavy flavour
- Higgs searches: Low-mass, SUSY



The DØ detector



Silicon tracker (SMT)

- <http://arxiv.org/abs/1005.0801>
- 6 barrels, 4 layers each, $z \sim 1$ m
+ **new Layer 0** @ $r = 1.6$ cm (RunIIb, see: <http://arxiv.org/abs/0911.2522>)
- Coverage $|\eta| < 2.5$

Central Fiber Tracker (CFT)

- 8 layers of scintillating fiber (axial and stereo)
- $20 < r < 51$ cm in **2T magnetic field**

Muon system covers $|\eta| < 2$

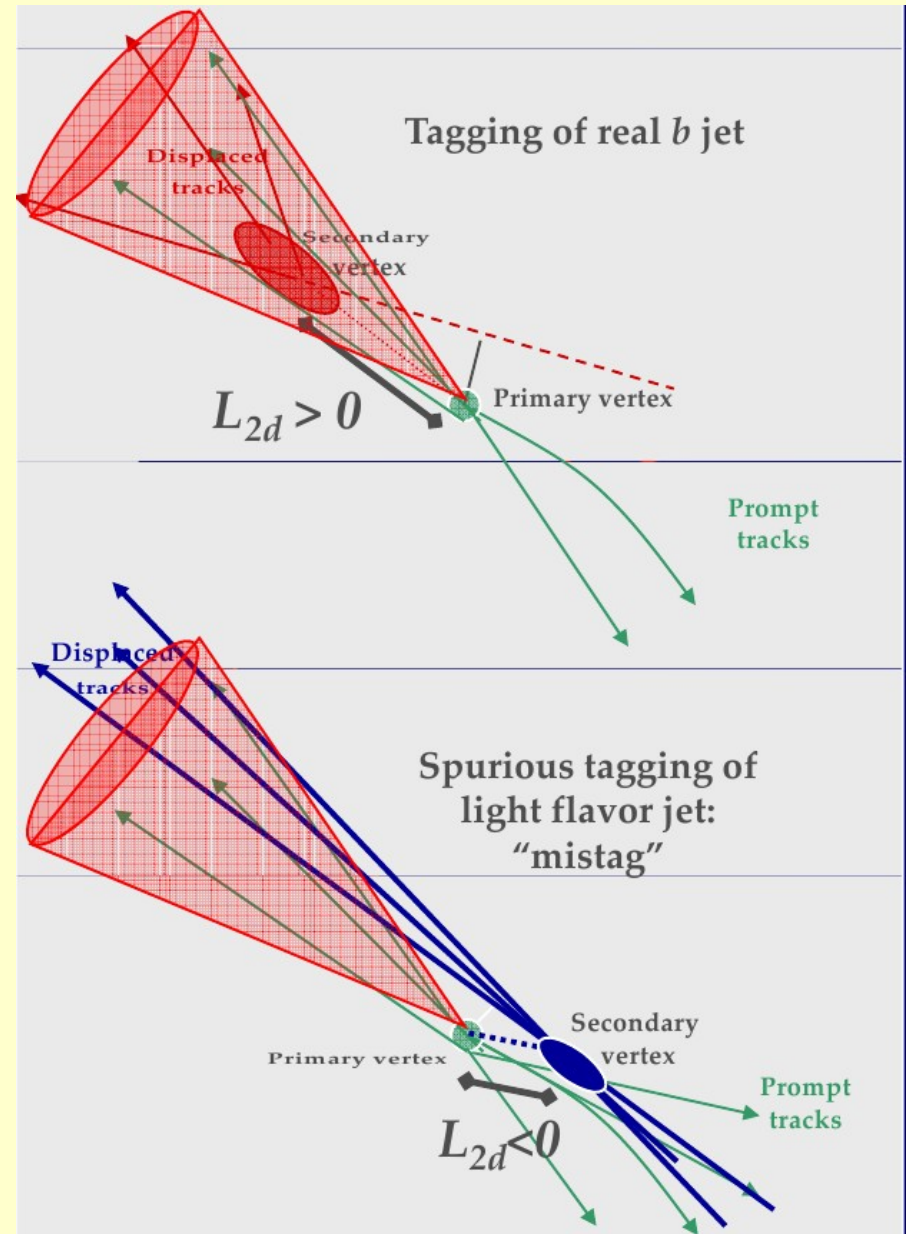
Tag vs. mis-tags

B hadrons properties

- Mass: $\sim 5 \text{ GeV}/c^2$
- Decay length: $\sim 3 \text{ mm}$
- **Hard** fragmentation
- Semi-leptonic decays

Fake / Mis -tags

- Primary vertex resolution
- Track parameters resolutions
- Long lived particles
- Secondary interactions



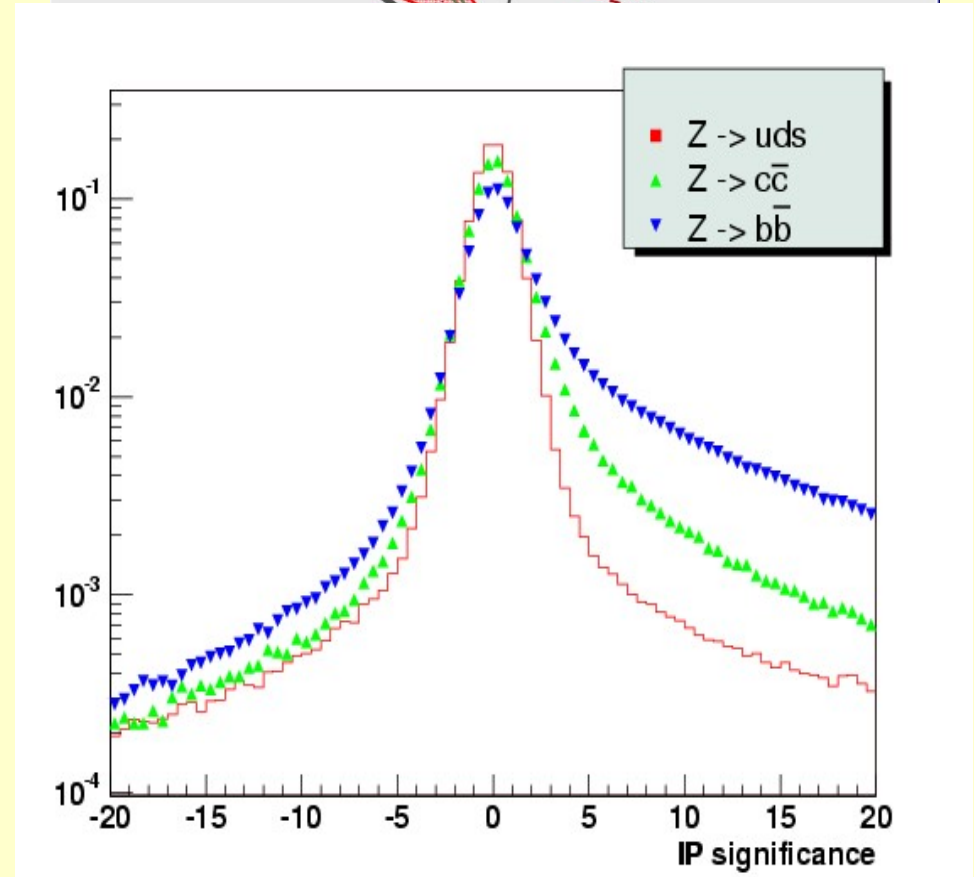
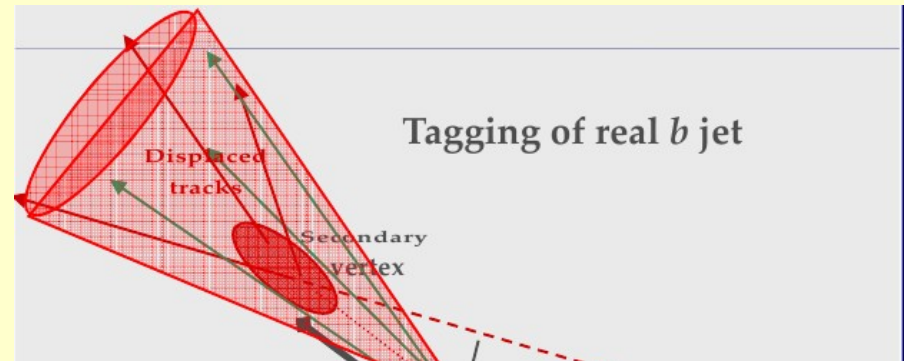
Tag vs. mis-tags

B hadrons properties

- Mass: $\sim 5 \text{ GeV}/c^2$
- Decay length: $\sim 3 \text{ mm}$
- **Hard** fragmentation
- Semi-leptonic decays

Fake / Mis -tags

- Primary vertex resolution
- Track parameters' resolution
- Long lived particles
- Secondary interactions

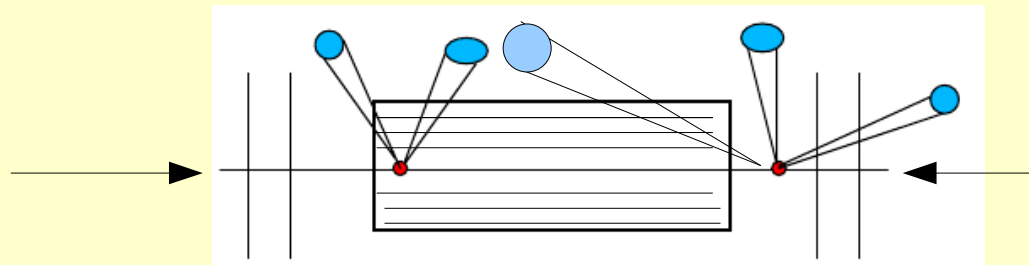


Tagging prerequisites

Taggability

In the following tagging algorithms are only based on **tracking and vertexing** of charged particles

- only (calorimeter) jets with minimum tracking information are considered
 - Interaction region, $\sigma_z \approx 25\text{cm}$, + detector acceptance affect track reconstruction efficiencies
 - performance dependence on η and interaction point's Z coordinate.

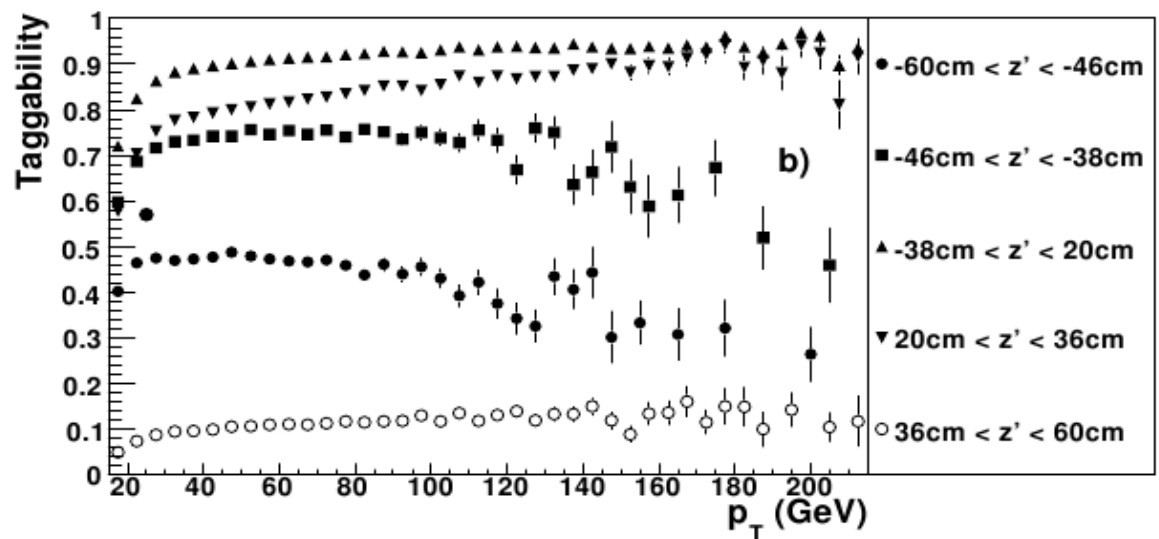
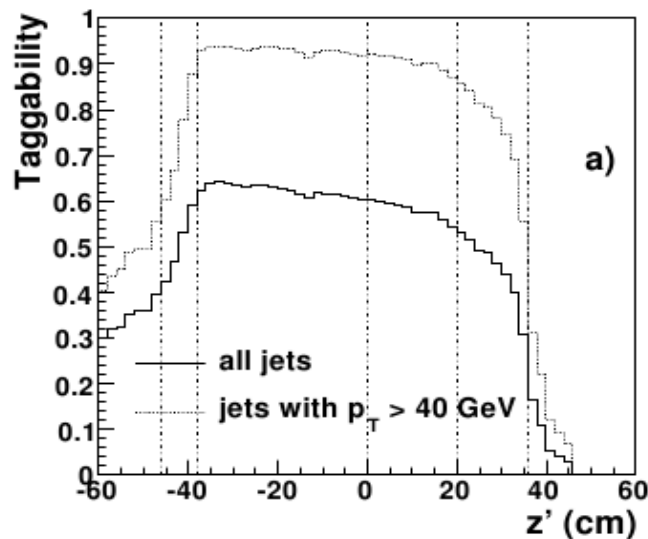


- Fraction of fake jets is ~small, but depends on **final** state.
 - Decoupling this effect from the tagging algorithms properly allows the extraction of a tagging performance which can be assumed to be **universal**, i.e., applicable to *any general final states*

Tagging prerequisites (II)

Taggable jets are thus defined as follow:

- **2-step clustering:**
 - i. along beam axis ($dca_z < 4 \text{ mm}$)
 - ii. 0.5 cone (snow mass) jets (*within each z-cluster*)
 ➡ finally require: $\Delta R(\text{calo-jet}, \text{track-jet}) < 0.5$
- **Track-jets:** 1 SMT hit tracks, seed track $p_T > 1 \text{ GeV}/c$, $p_T > 0.5 \text{ GeV}/c$ for other tracks
- **Parametrized as:** $F(p_T, \eta, z')$, with $z' \equiv |z| \cdot \text{sign}(\eta \cdot z)$



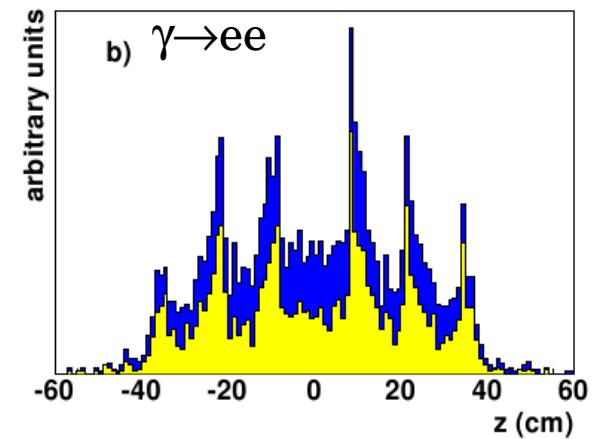
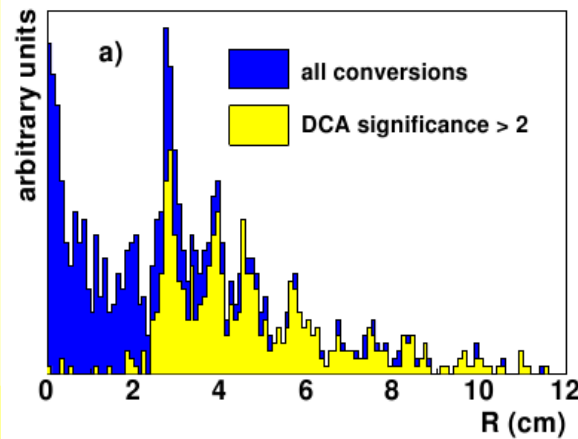
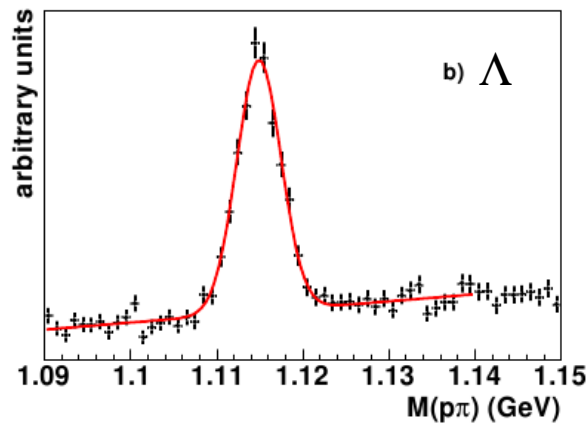
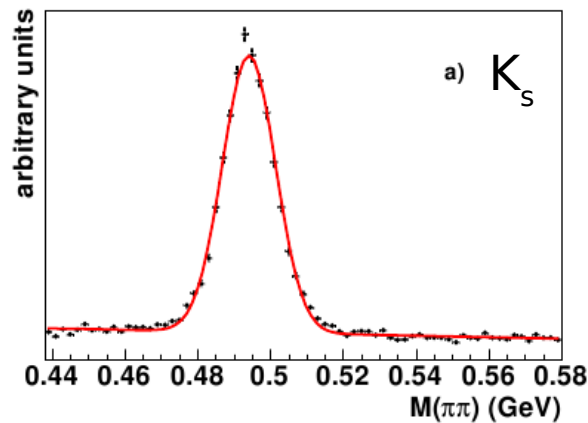
Tagging prerequisites (III)

Track preselection

- Each b-id. algorithm uses its own track reconstruction quality criteria

V^0 removal

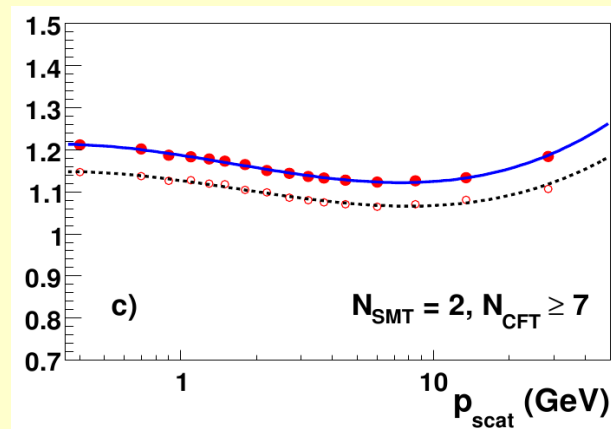
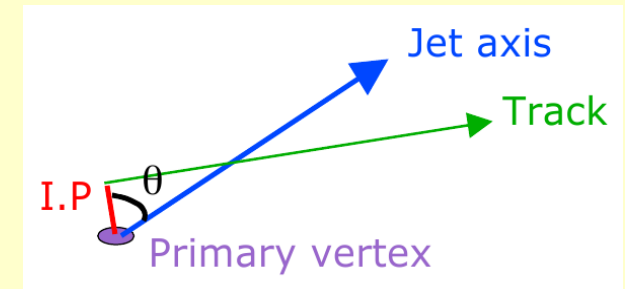
- Light strange hadrons have long lifetimes
- Photon conversions can occur at large distances in the tracker material



Algorithms (I)

Impact Parameter (IP) based tagger

- IP and its significance S_{IP} are **signed** w.r.t jet direction
- IP error calibrated in data and simulation for *multiple-scattering effects* and *PV resolution dependence*



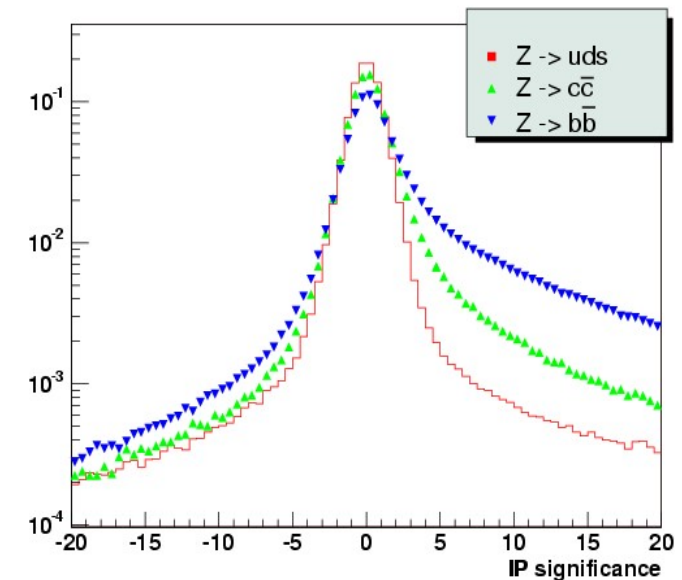
$$p_{\text{scat}} = p \sin\theta^{3/2}$$

Discrete (CSIP)

- counts tracks with: $|S_{IP}| > \text{cut}$ ($2 > 3$ | $3 > 2$)

Continuous (JLIP)

- p.d.f from negative IP resolution function, $R(s)$



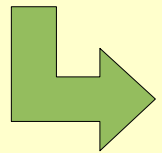
Algorithms (II)

Impact Parameter (IP) based tagger

Continuous (JLIP)

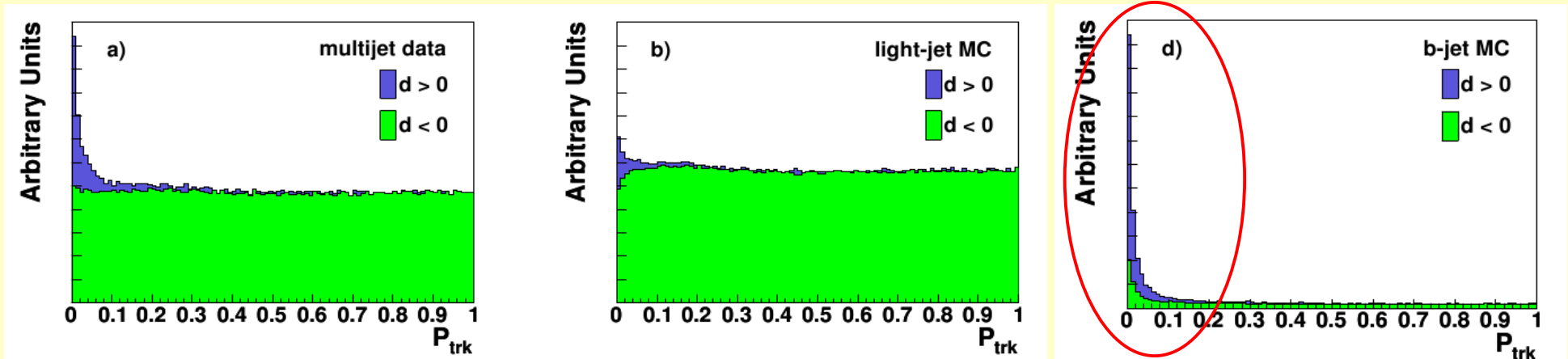
- p.d.f from negative IP resolution function, $\mathbf{R(s)}$

$$\mathcal{P}_{\text{trk}}(\mathcal{S}_d^{\text{corr}}) = \frac{\int_{-\infty}^{-|\mathcal{S}_d^{\text{corr}}|} \mathcal{R}(s) ds}{\int_{-\infty}^0 \mathcal{R}(s) ds}$$



$$\mathcal{P}_{\text{jet}}^{\pm} = \Pi^{\pm} \times \sum_{j=0}^{N_{\text{trk}}^{\pm}-1} \frac{(-\log \Pi^{\pm})^j}{j!} \quad \text{with} \quad \Pi^{\pm} = \prod_{i=1}^{N_{\text{trk}}^{\pm}} \mathcal{P}_{\text{trk}}(\mathcal{S}_{IP<0}^{IP>0})$$

Note: one can consider any set of tracks and e.g build an “*hemisphere-probability*” ($Z \rightarrow b\bar{b}$, LEP)



Algorithms (II)

Impact Parameter (IP) based tagger

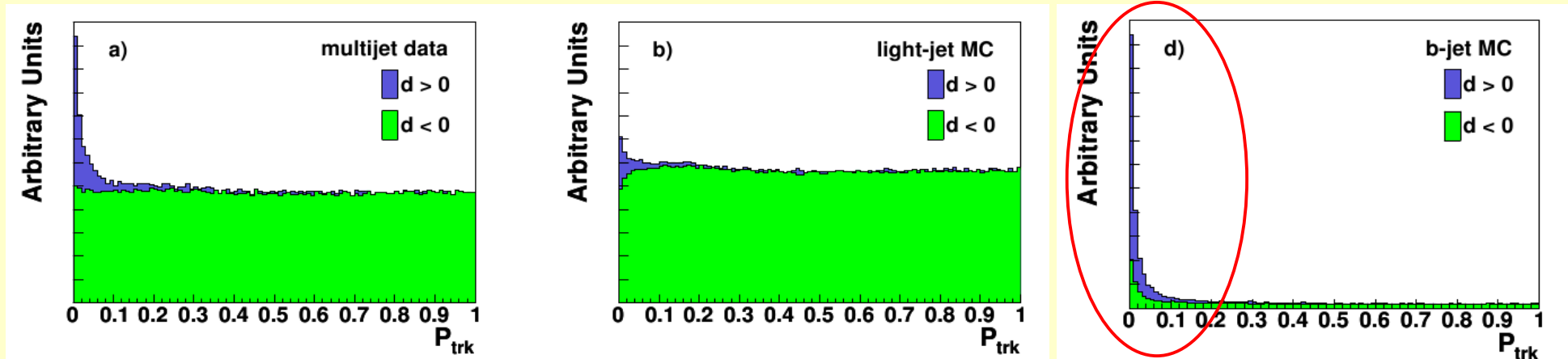
Continuous (JLIP)

- p.d.f from negative IP resolution function, $\mathbf{R(s)}$

$$\mathcal{P}_{\text{trk}}(S_d^{\text{corr}}) = \frac{\int_{-\infty}^{-|S_d^{\text{corr}}|} \mathcal{R}(s) ds}{\int_{-\infty}^0 \mathcal{R}(s) ds}$$

**For TMVA aficionados,
this is what is called “Rarity” in
the TMVAGui.C**

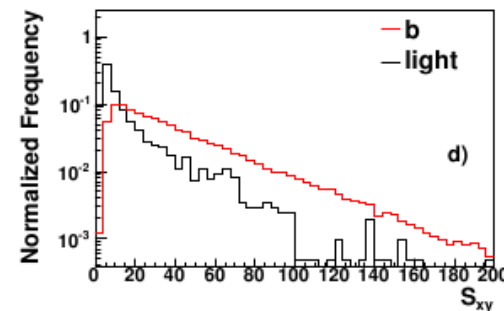
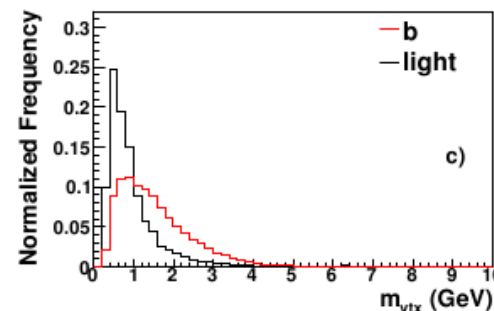
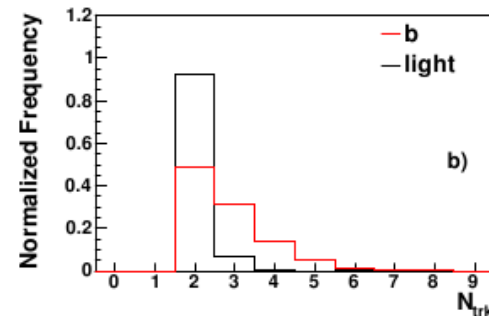
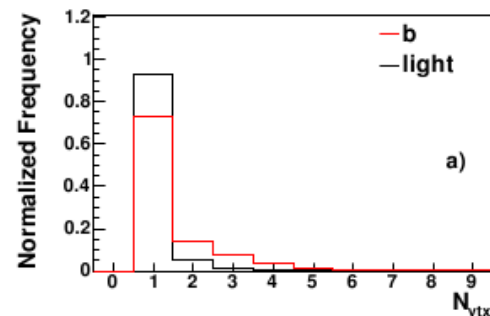
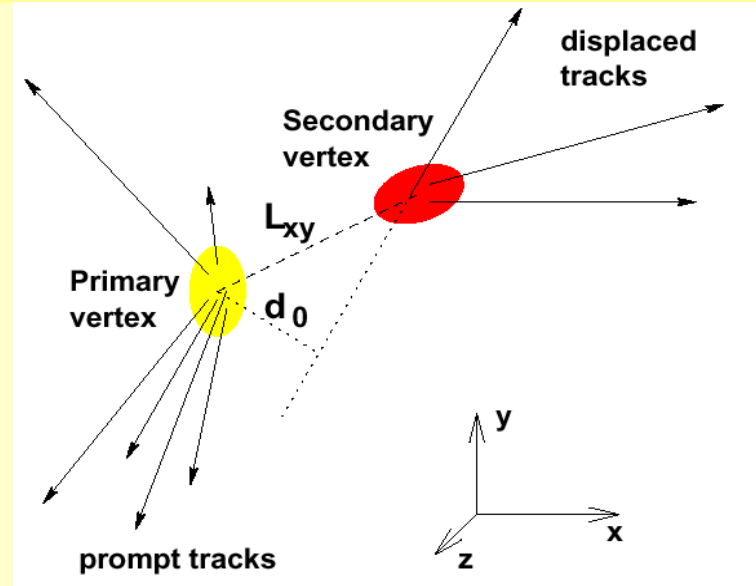
Note: one can consider any set of tracks and e.g build an “*hemisphere-probability*” ($Z \rightarrow b\bar{b}$, LEP)



Algorithms (III)

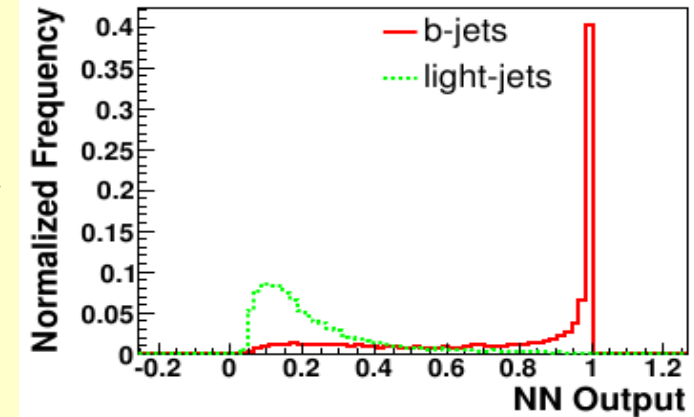
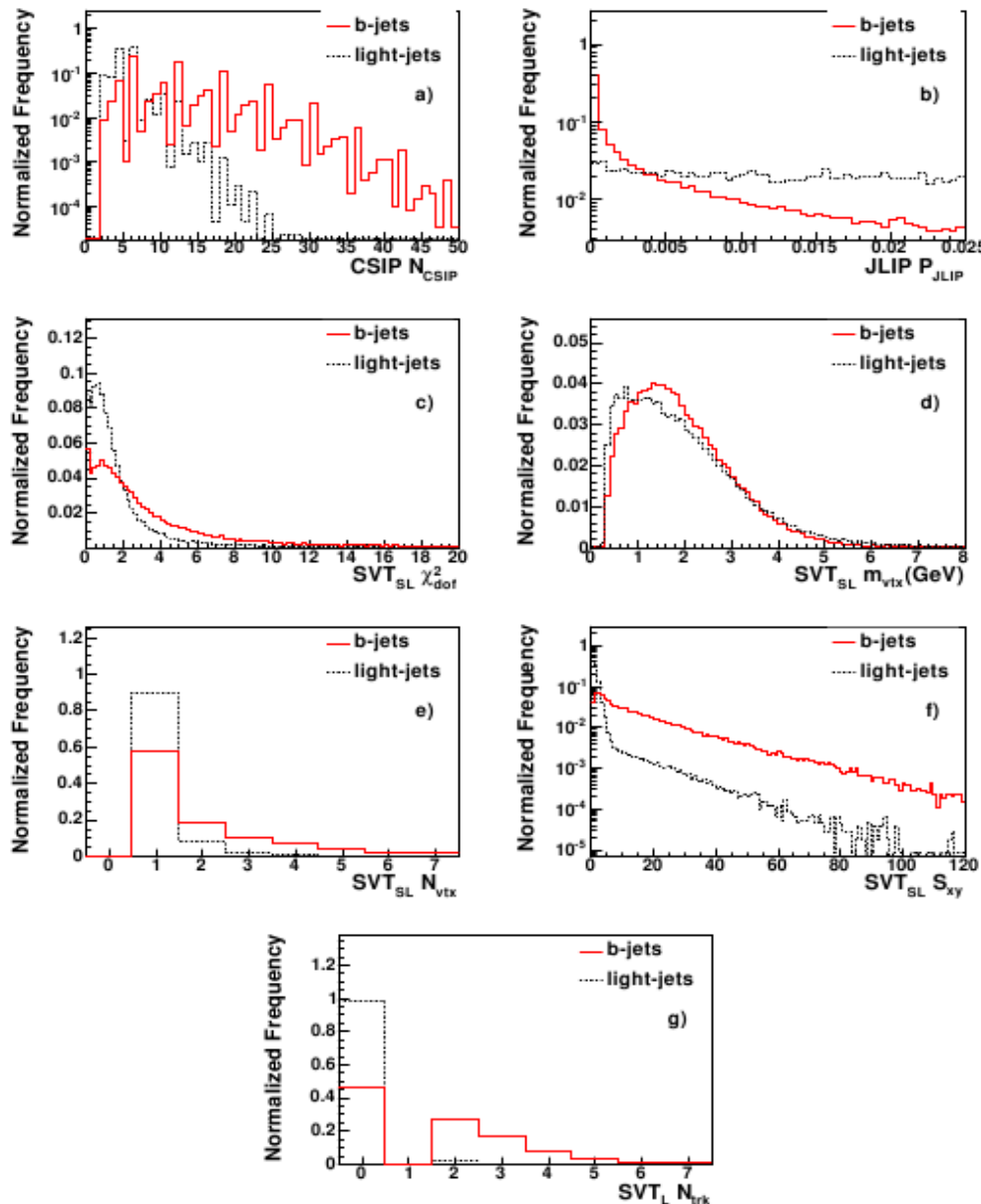
Secondary vertex, SVT

- Starts from track-based jets (*simple cone algo.*)
- Kalman-filter based vertex finder
- Track pruning w.r.t χ^2 contribution to vertex
- Tag is defined if:**
 $\Delta R(\text{vertex}, \text{jet}) < 0.5$ and if
decay length significance, $S_{Lxv} > \text{cut}$

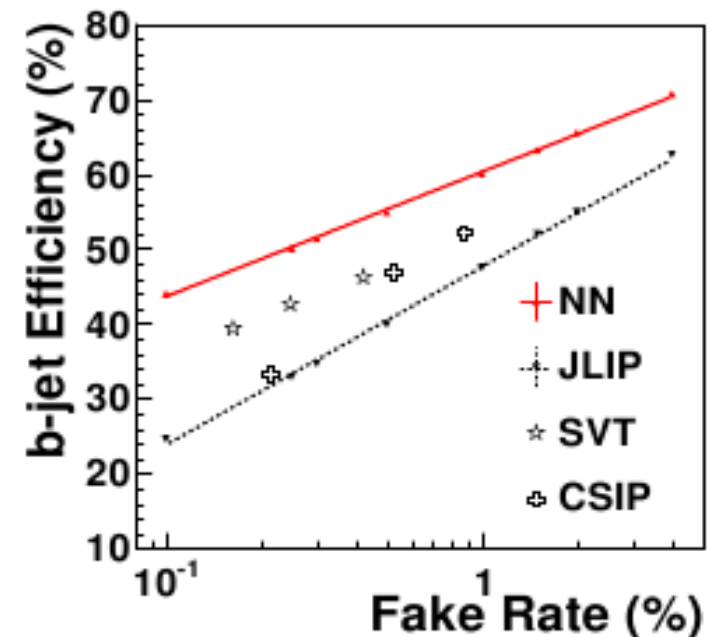


All in one: Neural Network tagger

Optimized selection of inputs: CSIP, JLIP & 5 SVT properties



... can lead to significant improvement:

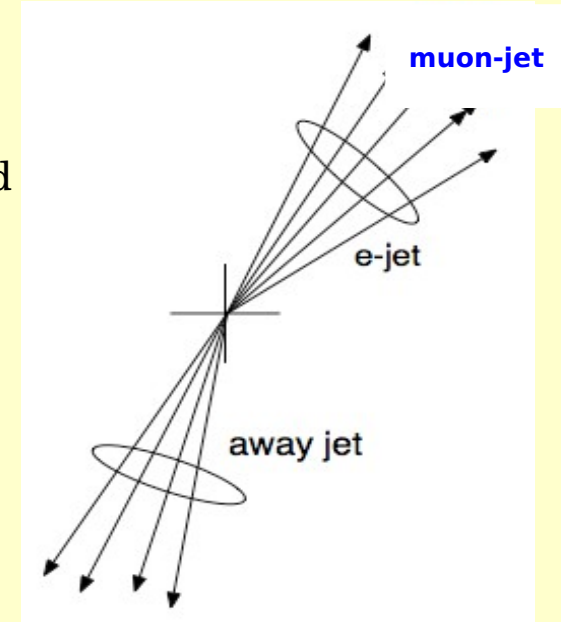
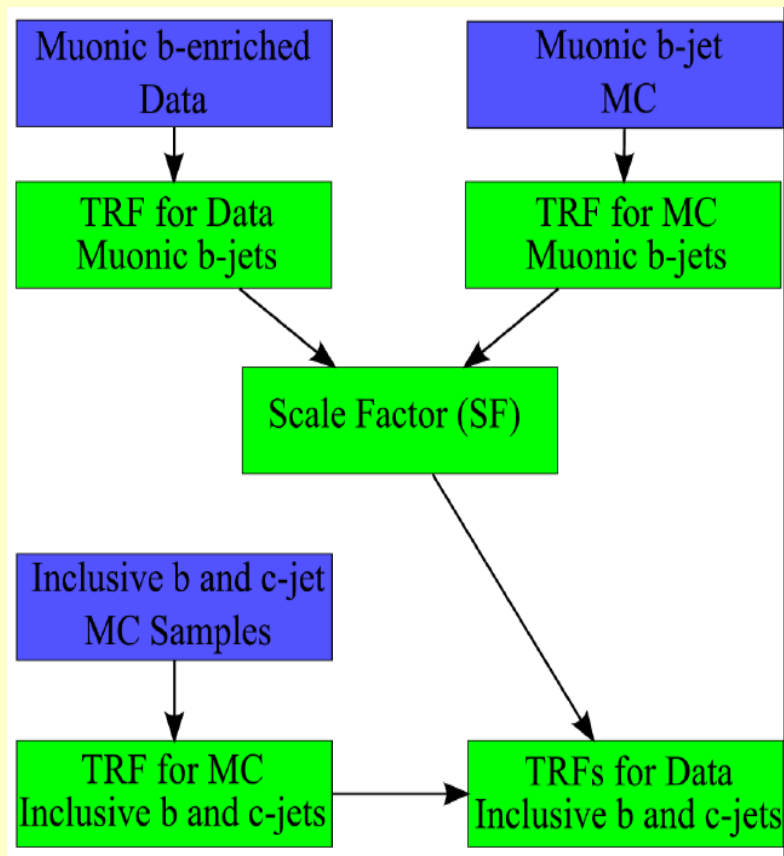


Performance measurements

B-identification efficiency (I)

Measured in data

- Using b-enriched data samples:
Di-jet back-to-back sample & require $\Delta R(<0.5)$ matched soft ($> 4\text{GeV}/c$) muon in jet
- Efficiency extracted using **SystemD** method



- Muonic data/MC b-Scale Factor:
 $\text{SF}_{b \rightarrow \mu}(\mathbf{p}_T, \eta)$
- Apply SF to **inclusive** b & c Tag Rate Function (TRF):

$$\epsilon_b^{\text{data}} = \frac{\epsilon_{b \rightarrow \mu X}^{\text{data}} \cdot \epsilon_b^{\text{MC}}}{\epsilon_{b \rightarrow \mu X}^{\text{MC}}} = \text{SF}_b \cdot \epsilon_b^{\text{MC}}$$

B-identification efficiency (II)

The SystemD method

- Historically developed to measure efficiency *solely in data*
- Simulation only used for **corrections factors** (MC/MC ratios)
- Main idea:** use *uncorrelated* selection criteria (i.e taggers) applied to various data samples and build a system of (non-linear) equations

General case:

- Consider **s** data samples composed of 1 signal and **f** backgrounds. Each sample **j** can gives **1+f unknowns**: the signal and backgrounds fractions :

$$n_{i=0\dots f}^{j=1\dots s} \text{ constrained by: } \sum_{i=0}^f n_i^j = 1$$

- Each tagger **k** gives also 1+f unknowns, the **efficiencies** : $\varepsilon_{i=0\dots f}^{k=1\dots t}$
- When applying the tagger k on sample j, only a fraction **q_j^k** of the total number of events survives:

$$q_j^k = \sum_{i=0}^f \varepsilon_i^k n_i^j$$

- When applying e.g 2 uncorrelated criteria:

$$q_j^{k_1, k_2} = \sum_{i=0}^f \varepsilon_i^{k_1} \varepsilon_i^{k_2} n_i^j$$

B-identification efficiency (III)

- Combining **t** taggers and **s** samples $\Rightarrow 2^t \cdot s$ equations
- To solve the system, one needs at least as many equations as unknowns:

$$2^t \cdot s \geq (1 + f)(s + t)$$

- The first combinations are:
 - $s = 2, t = 2, f = 1$: 8 equations and 8 unknowns;
 - $s = 1, t = 3, f = 1$: 8 equations and 8 unknowns;
 - $s = 2, t = 3, f = 2$: 16 equations and 15 unknowns;
 - $s = 6, t = 2, f = 2$: 24 equations and 24 unknowns.

In practice finding many samples and (uncorrelated) taggers is difficult

Note: $t = 1, s = 2, f = 1$ is known as the Matrix Method :-)

B-identification efficiency (IV)

SystemD and b-tagging

$s = 2$ and only the first combination is considered:

- **2 (uncorrelated) taggers:**

NN-tagger & soft lepton(muon) tagger w/ a p_T^{rel} cut

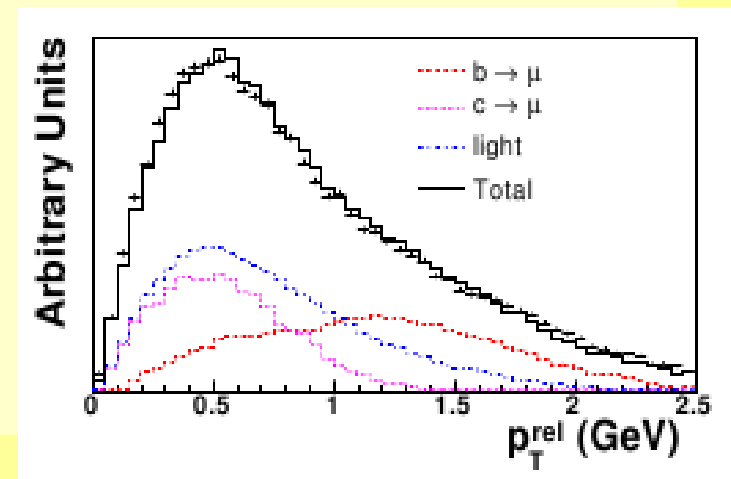
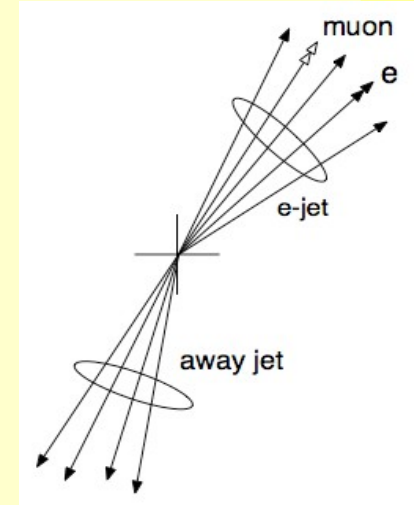
- **2 data samples** w/ different flavour content:

muon-jet & muon-jet + away tag

- Apply 2 taggers **separately / simultaneously** on 2 samples and solve (*analytically or numerically*) the 8 equations / 8 unknowns among which:

$$\epsilon_b(\text{NN})$$

- *c* and *light* jets are considered as a single background (*i.e* $f = 1$)



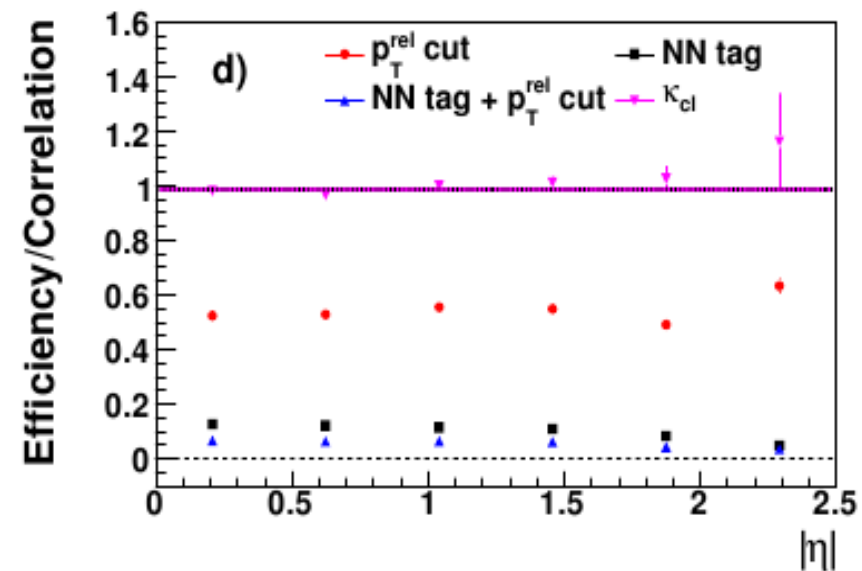
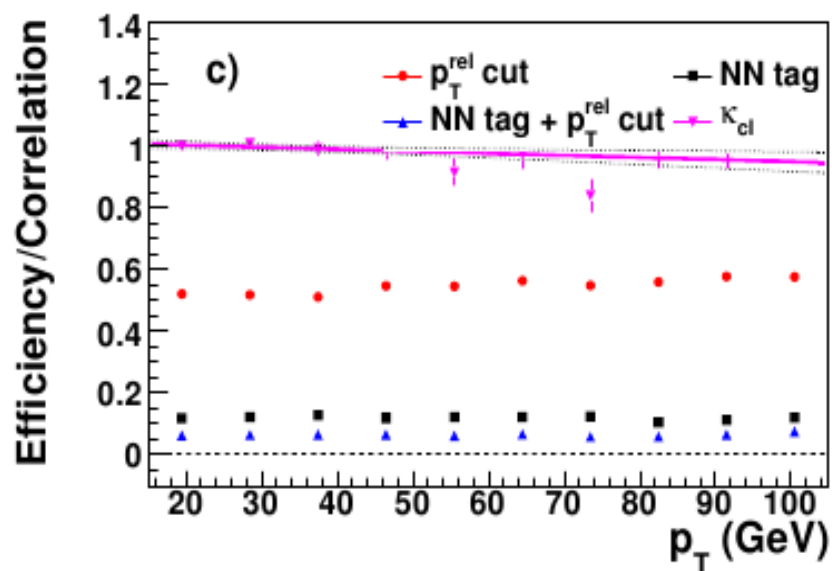
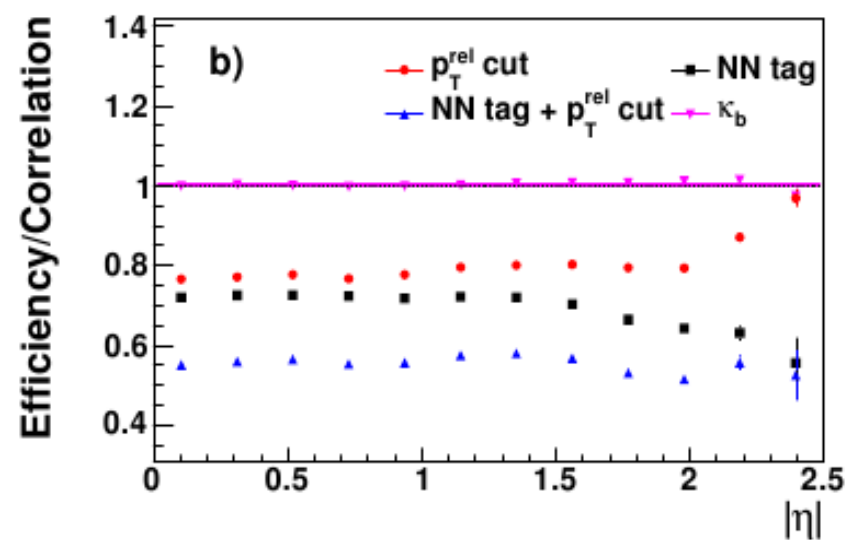
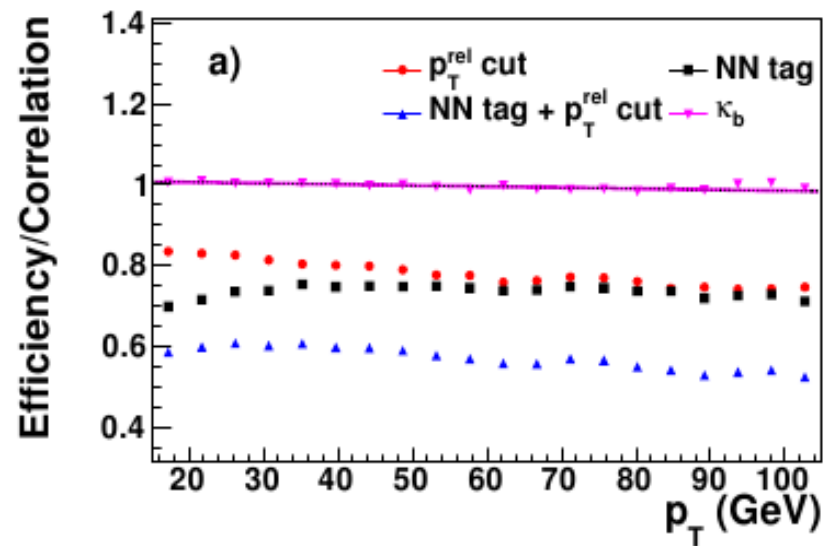
B-identification efficiency (V)

Corrections factors

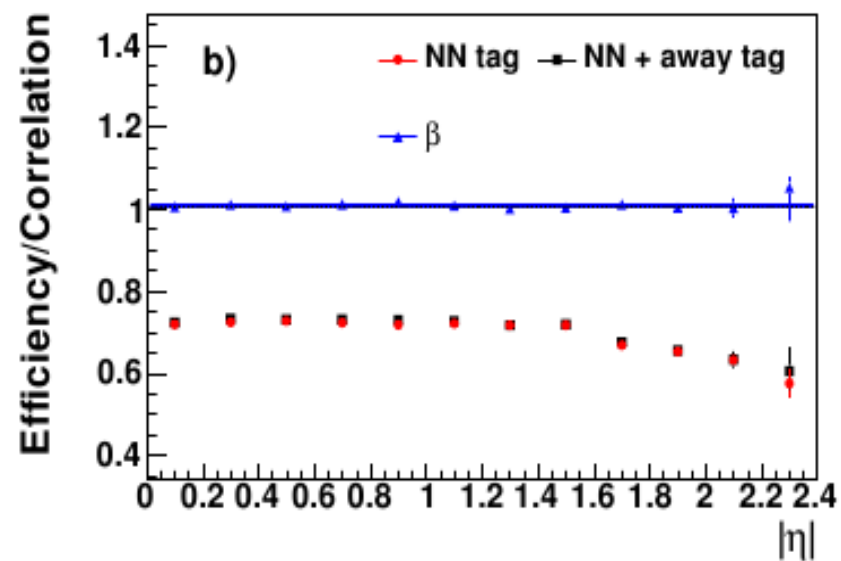
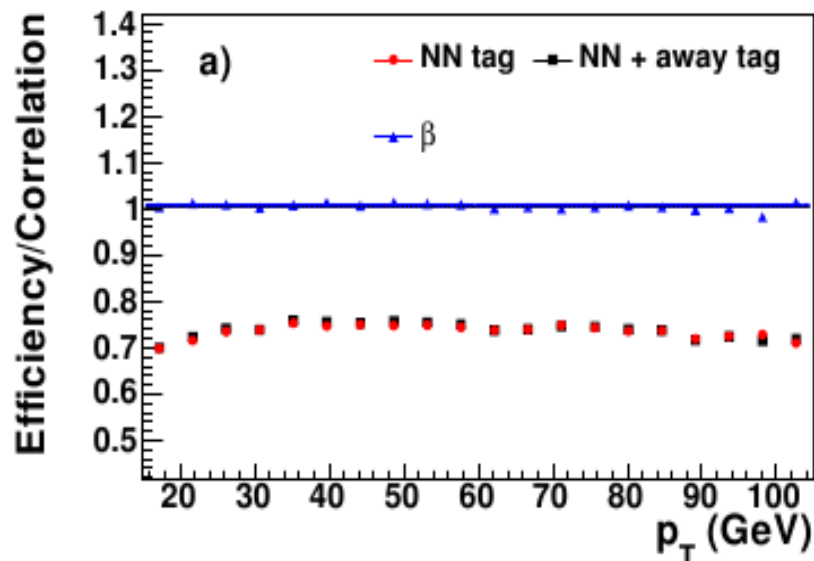
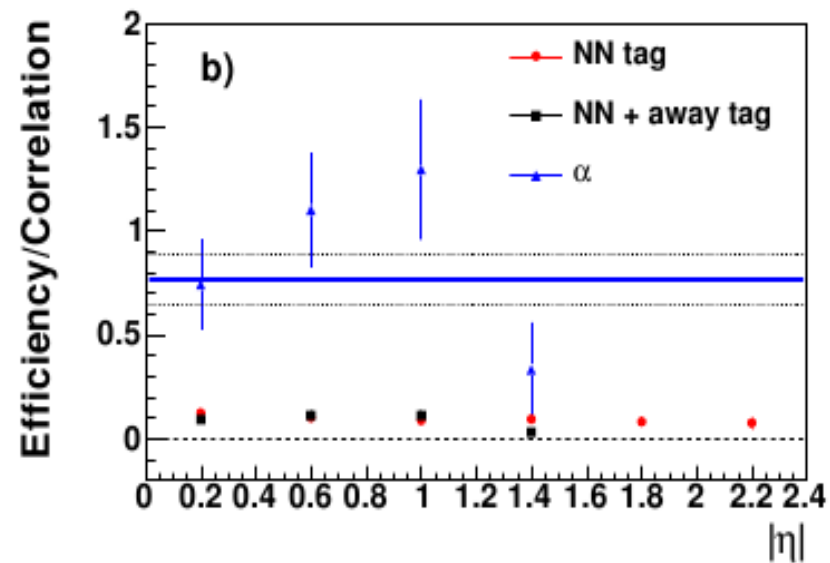
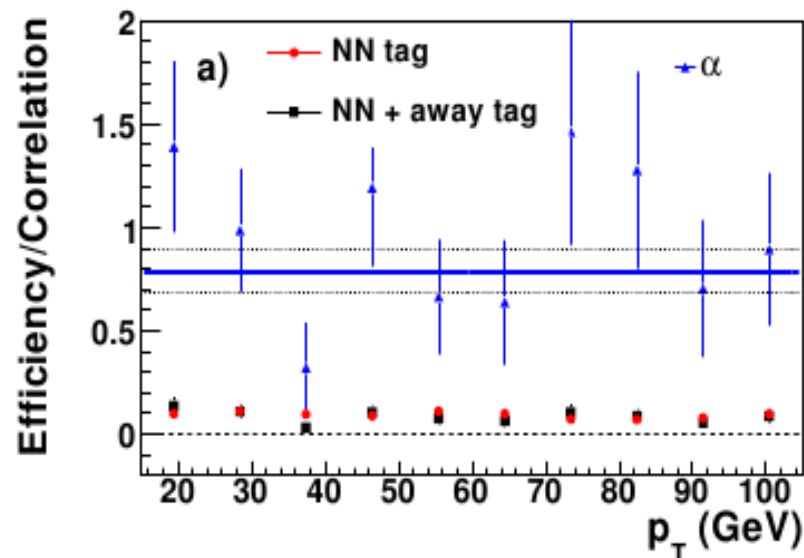
- The SLT and NN are assumed to be uncorrelated (*mass vs. lifetime*)
- The away-tag and the SLT are uncorrelated (*but same PV!*)
 - Introduce corrections factors for signal and backgrounds to quantify these correlations
 - Parameterized as a function of jet p_T /eta

f_b	+	f_{cl}	=	1
$f_{b\epsilon_b^l}$	+	$f_{cl\epsilon_{cl}^l}$	=	Q^l
$f_{b\epsilon_b^m}$	+	$f_{cl\epsilon_{cl}^m}$	=	Q^m
$f_{b\epsilon_b^b}$	+	$f_{cl\epsilon_{cl}^b}$	=	Q^b
$f_{b\kappa_b\epsilon_b^l\epsilon_b^m}$	+	$f_{cl\kappa_{cl}\epsilon_{cl}^l\epsilon_{cl}^m}$	=	$Q^{l,m}$
$f_{b\epsilon_b^m\epsilon_b^b}$	+	$f_{cl\epsilon_{cl}^m\epsilon_{cl}^b}$	=	$Q^{m,b}$
$f_{b\beta\epsilon_b^b\epsilon_b^l}$	+	$f_{cl\alpha\epsilon_{cl}^b\epsilon_{cl}^l}$	=	$Q^{b,l}$
$f_{b\kappa_b\beta\epsilon_b^l\epsilon_b^m\epsilon_b^b}$	+	$f_{cl\kappa_{cl}\alpha\epsilon_{cl}^l\epsilon_{cl}^m\epsilon_{cl}^b}$	=	$Q^{l,m,b}$

B-identification efficiency (VI)



B-identification efficiency (VII)



B-identification efficiency (VIII)

Systematic uncertainties

- Corrections factors measured with finite stat. MC
 - ↳ vary within errors 1 correction (*fix the others*) and re-run SystemD
- p_T^{rel} cut varied from 0.3 to 0.8 GeV/c
- Add all errors quadratically
- Apply in *each jet p_T and eta* bins for *each operating point (OP)*

SystemD syst. errors:

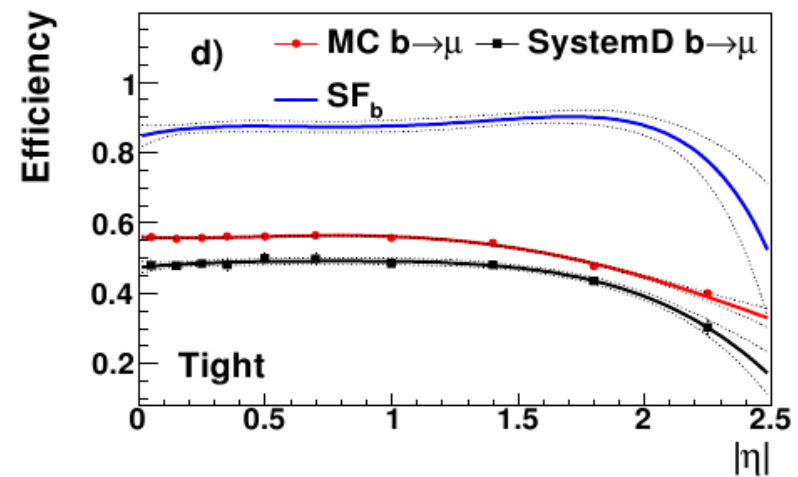
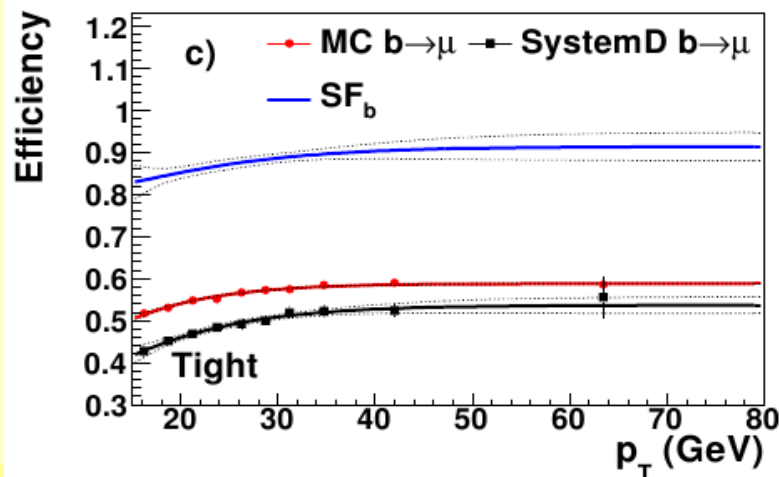
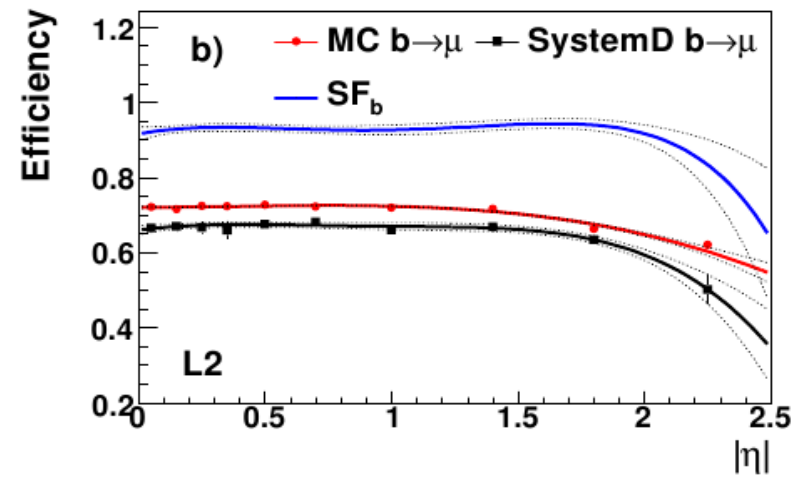
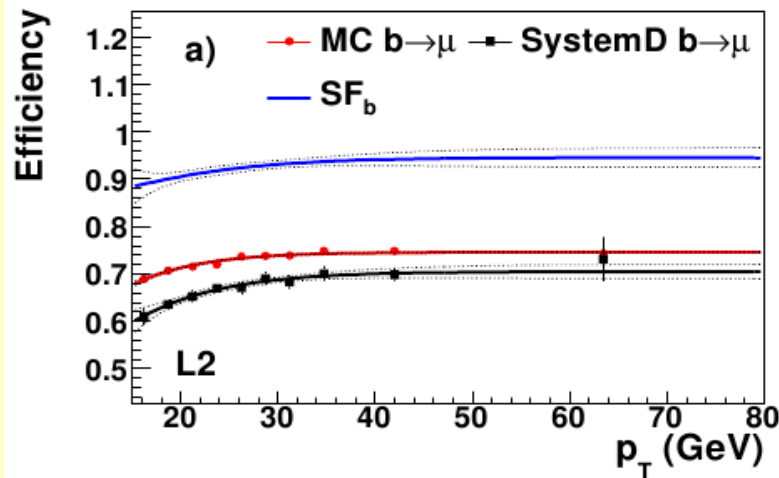
	L2	Tight
Efficiency	65.9%	47.6%
α	0.0%	0.0%
β	0.2%	0.6%
κ_b	0.7%	1.2%
κ_{cl}	0.3%	0.2%
p_T^{rel}	1.0%	0.7%
<i>SystemD</i> Total	1.3%	1.5%

- **B-jet efficiencies errors: ~2% to ~5%**

B-identification efficiency (IX)

Scale factors are measured for 12 operating points

- Optimize efficiency / purity depending on physics channels
e.g single / double (*asymmetric*) tags, ...

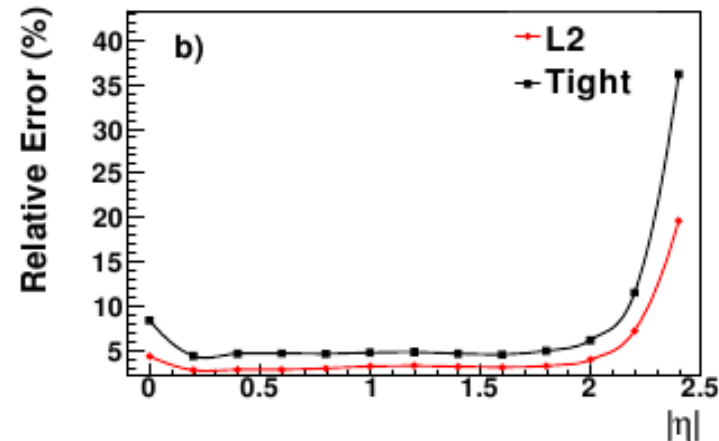
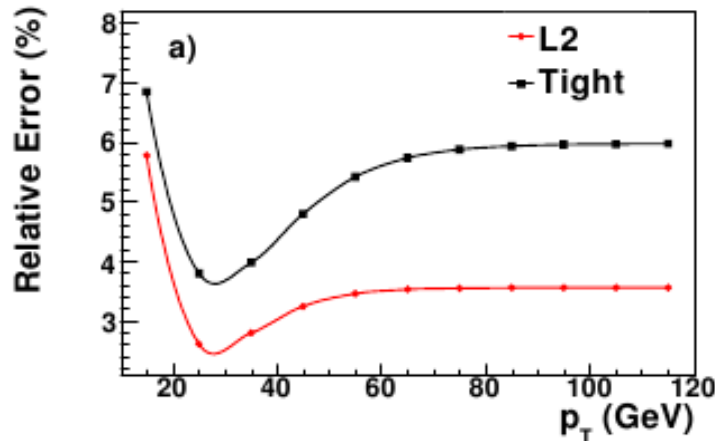


B-identification efficiency (IX)

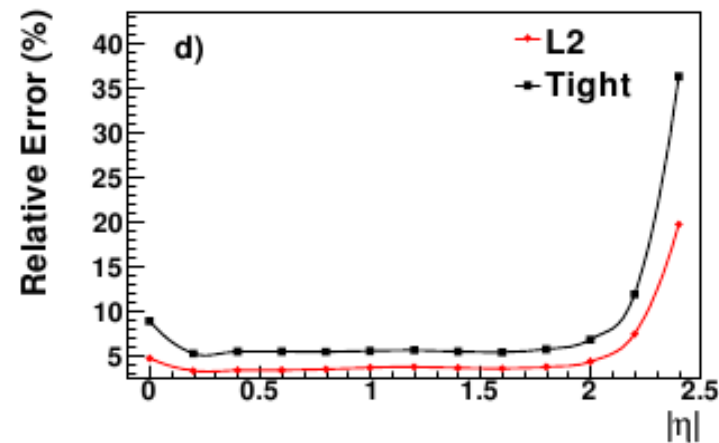
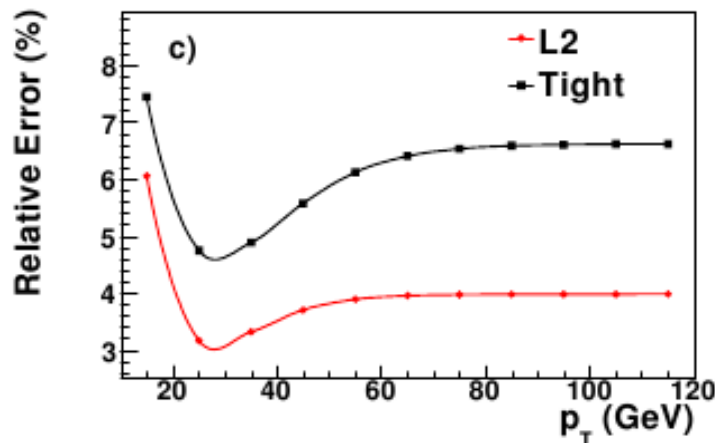
Scale factors are measured for 12 operating points

- Optimize efficiency / purity depending on physics channels
e.g single / double (*asymmetric*) tags, ...

SF_b



TRF_b



Fake rate

Goal

- Estimate $\varepsilon_{\text{light}}$ where light = u,d,s and gluon
- Measured in *data*

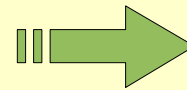
Estimated from *negative tags*

Corrected for:

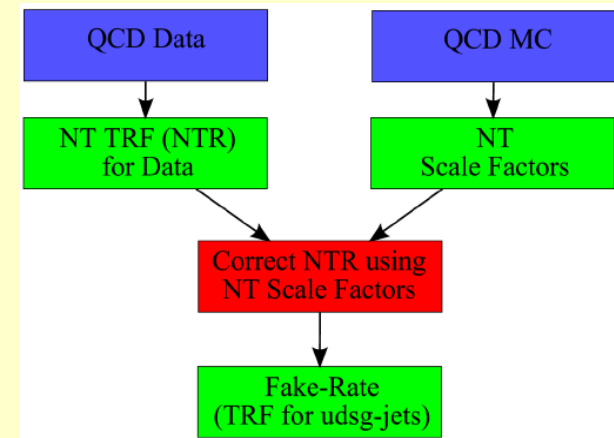
- HF contamination:
- Neg./Pos. asymmetry:

$$F_{\text{hf}} = \varepsilon_{\text{QCD,light}}^- / \varepsilon_{\text{QCD,all}}^-$$

$$F_{\text{lf}} = \varepsilon_{\text{QCD,light}}^+ / \varepsilon_{\text{QCD,light}}^-$$



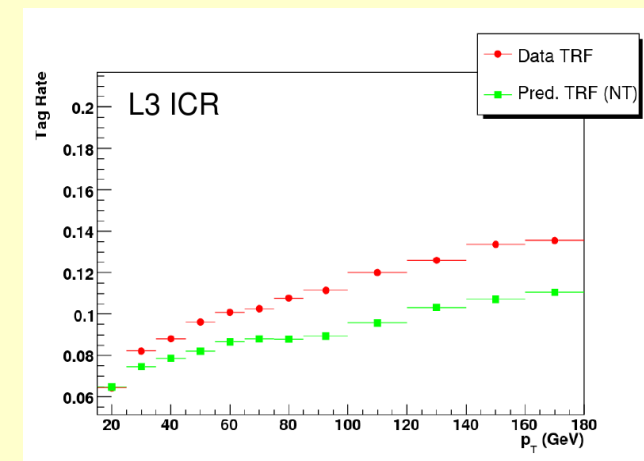
$$\varepsilon_{\text{light}} = \varepsilon_{\text{data}}^- \cdot F_{\text{hf}} \cdot F_{\text{lf}}$$



Parameterisation

- $F(p_T, \eta(\text{CC, ICR, EC}))$

But: NT method **underestimates** fake-rate (*hidden* in “experimental” *k*-factor)



Fake rate

Goal

- Estimate ϵ_{light} where light = u,d,s and gluon
- Measured in *data*

Estimated from *negative tags*

Corrected for:

- HF contamination:
- Neg./Pos. asymmetry:

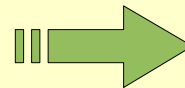
Parameterisation

- $F(p_T, \eta(\text{CC}, \text{ICF}))$

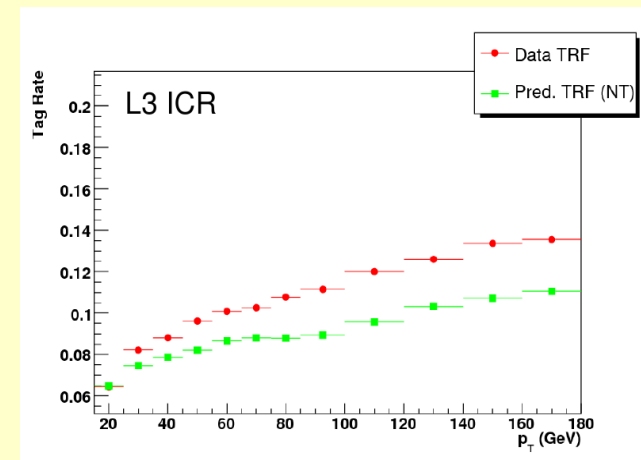
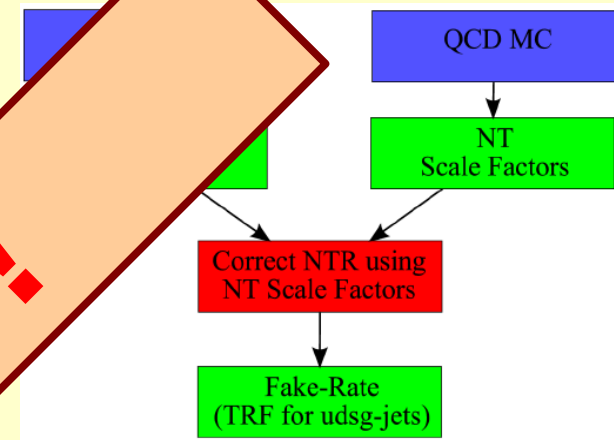
But: NT method **underestimates** fake-rate (*hidden* in “experimental factor”)

$$F_{\text{hf}} = \epsilon_{\text{Q}}^-$$

$$F_{\text{lf}} = \epsilon_{\text{CD,light}}^-$$



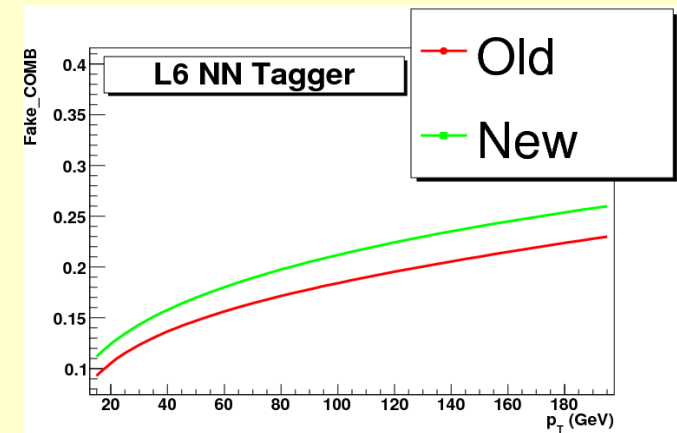
$$\epsilon_{\text{light}} = \epsilon_{\text{data}}^- \cdot F_{\text{hf}} \cdot F_{\text{lf}}$$



Fake rate

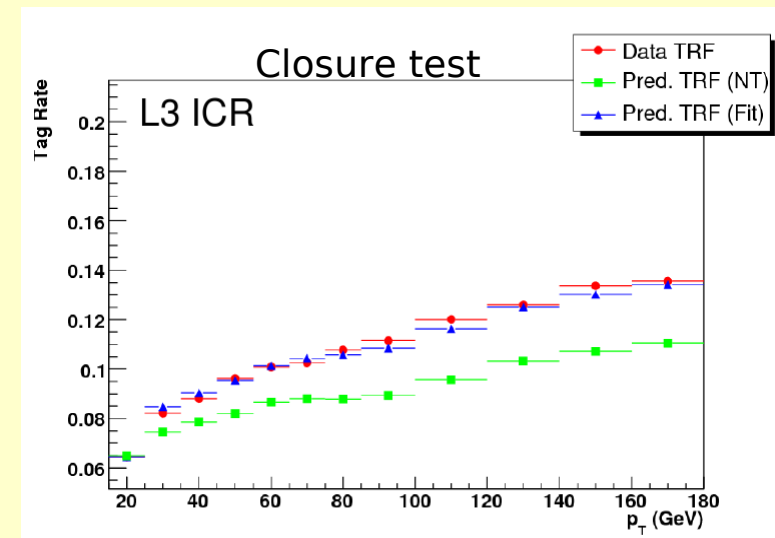
New method (default since summer 2009)

- More data-driven
- Can be applied on specific dataset (e.g hbb)
- Uses b/c tag rate from System8



Results

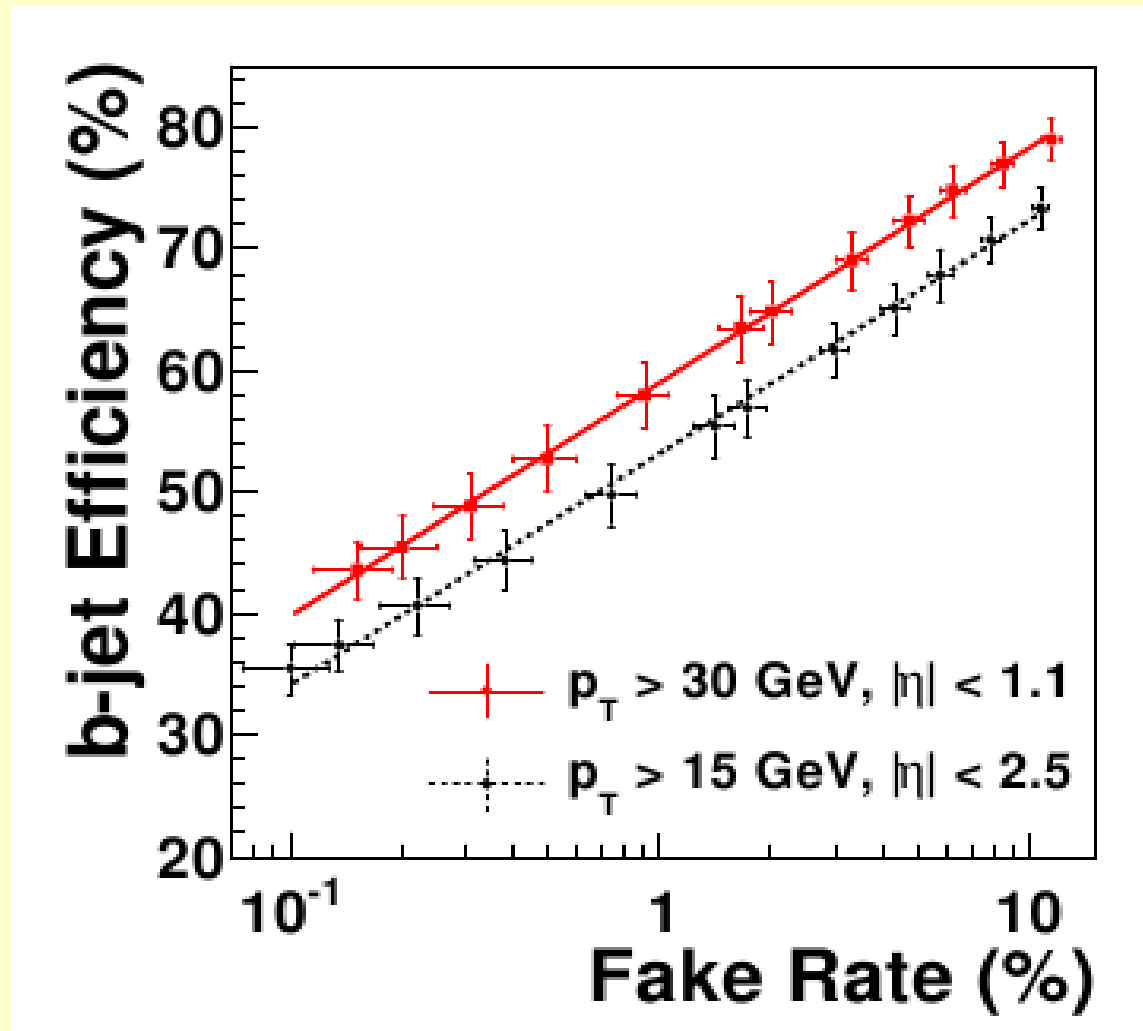
- Similar shapes
- 20-50% higher rate
- Good closure tests.
- $K_{hf}^{exp} \sim 1$



Performance

Final data performance

- Using MC Z decays with data / MC scale factors



Conclusion

- These algorithms have been used in many publications of D0RunII analyses
- SystemD applied to b-identification efficiency measurement is a powerful method with **little dependency on simulation**
 - *It is already used in both ATLAS and CMS*
- **Other / On-going / Future developments:**
 - MVA taggers
 - Improvements to the algorithms and methods
 - Fake track killer / tracking tuning / neg. tags / ...
 - See talks about Tevatron run extension ... :(
 - Better understanding of detector response
 - But also higher instantaneous luminosity
 - *You can contribute !*

Thank You !

Back-up

System D in simulated events

Method validation in simulation:

